Outline

- 1. Launching a Spark cluster on AWS
- 2. HDFS, YARN and Spark
- 3. First steps with Spark
- 4. Map-and-reduce architecture

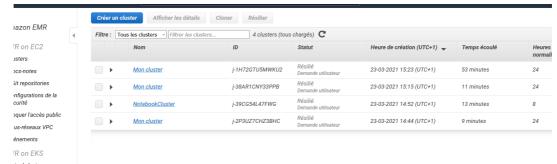
Création d'un cluster Spark sur AWS

Une fois connecté à la console AWS cherchez le service EMR pour *Elastic Map Reduce*. C'est la *Plateforme as a Service* qui permet de créer un cluster Hadoop dans AWS. Vous allez juste choisir la configuration de votre cluster, et AWS va créer toutes les VM, les mettre en réseau et installer toutes les applications choisie pour vous. Créer un cluster Hadoop à la main est laborieux et n'est pas réellement intéressant, c'est pourquoi les divers fournisseurs de services cloud proposent de telles PaaS

☐ Vous allez arriver sur un écran similaire à celui-ci



Les fois suivantes il ressemblera à cela



Dans tous les cas cliquez sur Créer un cluster

☐ Sur la page suivante voici les configurations à saisir :

oxdot Nom du cluster : le nom	que vous	souhaitez.	Le nom p	oar défaut	convient t	:rès
bien						

☐ Journalisation : décochez ce paramètre. Il permet de sauvegarder les *log* (journaux) de votre cluster pour ensuite allez chercher la source d'une erreur. Mais pour ce TP cela ne sera pas utile

Mode de lancement : Cluster . La différence entre Cluster et Exécution d'étape est que Cluster permet interagir avec le cluster, alors que qu' Exécution d'étape créer un cluster, réalise lance un ou plusieurs scripts et s'arrête. C'est parfait quand vous voulez lancer un "job" sur vos données. Le cluster se lance, fait les calculs puis va exporter les résultats et s'éteindre tout seul.						
	Libérée (Release en VO) : emr-5.31.0, donc l'avant dernière version 5.XX. La dernière à des problèmes pour l'utilisation de notebook.					
Type d'instance : des m5.xlarge conviennent parfaitement. Si vous voulez vous pouvez essayer des machines plus puissantes, mais cela ne va pas impacter fortement les temps de calculs.						
Nombre d'instance : ③, mais vous pouvez essayer avec plus d'instance (limitez vous à 6).						
	r vous donner une idée des prix unitaire des machines voici un tableau des ance m5.XX. Pour le prix du cluster, multipliez par le nombre de machine					
	•					
	•					
instance m5.XX. Pour le pri	x du cluster, multipliez par le nombre de machine					
instance m5.XX. Pour le pri	x du cluster, multipliez par le nombre de machine Prix unitaire par heure					
instance m5.XX. Pour le pri Instance m5.xlarge	Prix unitaire par heure 0.24\$					
instance m5.XX. Pour le pri Instance m5.xlarge m5.2xlarge	Prix unitaire par heure 0.24\$ 0.48\$					
instance m5.XX. Pour le pri Instance m5.xlarge m5.2xlarge m5.4xlarge m5.8xlarge	Prix unitaire par heure 0.24\$ 0.48\$ 0.96\$					
instance m5.XX. Pour le pri Instance m5.xlarge m5.2xlarge m5.4xlarge m5.8xlarge □ Paire de clé EC2 : sélections	Prix unitaire par heure 0.24\$ 0.48\$ 0.96\$ 1.86\$ nez la clé du TP 0. Si vous n'en avez pas sélectionnez					

Créer un cluster - Options rapides Accéder aux options avancées Configuration générale Nom du cluster Mon cluster Journalisation 0 Mode de lancement Cluster Exécution d'étape Configuration des logiciels Libére Core Hadoop: Hadoop 2.10.0, Hive 2.3.7, Hue 4.7.1. HBase: HBase 1.4.13, Hadoop 2.10.0, Hive 2.3.7, Hue 4.7.1, Phoenix 4.14.3, and ZooKeeper 3.4.14 Presto: Presto 0.238.3 with Hadoon 2.10.0 HDFS Spark: Spark 2.4.6 on Hadoop 2.10.0 YARN and Zeppelin 0.8.2 Utiliser AWS Glue Data Catalog pour les Configuration du matériel Le type d'instance sélectionné ajoute un volume EBS GP2 par défaut de 64 GiO par instance. En savoir plus 🛂 Type d'instance m5.xlarge Nombre d'instances 3 (1 nœud maître et 2 nœuds principaux) Cluster scaling Sécurité et accès Paire de clés EC2 big_data_ensai Apprenez à créer une paire de clés EC2. Autorisations Par défaut Personnalisé Utilisez les rôles IAM par défaut. Si des rôles sont absents, ils seront créés automatiquement pour vous avec des stratégies gérées pour les mises à jour automatiques de stratégies. Rôle EMR EMR_DefaultRole 🖸 🕦 Profil d'instance EC2 EMR EC2 DefaultRole [7] € Annuler Créer un cluster La création d'un cluster prend plus de temps que la création d'une machine unique (de l'ordre de quelques minutes). Car AWS doit lancer X machines avec des configurations lourdes, les mettre en réseau etc. Une fois le cluster passer en "En attente" ou "Stand by", allez dans blocs-notes, puis cliquez sur Créer un bloc-notes



☐ Sur l'écran suivant vous allez devoir spécifier :

☐ Un nom pour votre <i>notebook</i>	
Le cluster que vous souhaitez utiliser. Noter que vous pouvez créer à la volée un cluster si vous le souhaitez	
Le rôle de sécurité et les groupes associés au service. Nous allons garder les paramètres par défaut	
L'endroit où sera stocké votre notebook sur S3. Cela peut servir à recharger u notebook fait plus tôt.	ın

☐ Un dépôt git si vous souhaitez versionner votre code.

Voilà quoi devrait ressembler votre écran avant validation.

Créer un bloc-notes

Nommer et configurer votre bloc-notes

Nommez votre bloc-notes Jupyter géré par EMR, choisissez un cluster ou en créez-en un et personnalisez les options de configuration si vous le souhaitez. En savoir plus 🛂 Nom du bloc-notes.* notebook Les noms peuvent contenir uniquement des lettres (a-z), des chiffres (0-9), des traits d'union (-) ou des traits de so Description Mon super notebook exemple Cluster* O Choisir un cluster existant Choose Mon cluster j-30YZFPZ86VI04 Créer un cluster 1 Groupes de sécurité

Utiliser des groupes de sécurité par défaut

Output

Unique de sécurité par défaut Choisir des groupes de sécurité (vpc-e2567a98) Rôle de service AWS* EMR_Notebooks_DefaultRole Emplacement du bloc- Choose an S3 location where files for this notebook are saved. Use the default S3 location s3://aws-emr-resources-523967347856-us-east-1/notebooks/ Choose an existing S3 location in us-east-1 ▶ Référentiel Git Lien vers un référentiel Git pour enregistrer votre bloc-notes dans un environnement de gestion des ve sions ▶ Balises ① * Obligatoire Créer un bloc-notes

La création du *notebook* doit être rapide. Une fois votre *notebook* prêt cliquez sur ouvrir dans JupyterLab. Cela ouvrira une interface JupyterLab pour rédiger des *notebooks*. Par défaut vous pouvez faire des *notebooks*:

☐ Python3

PySpark: Spark avec python

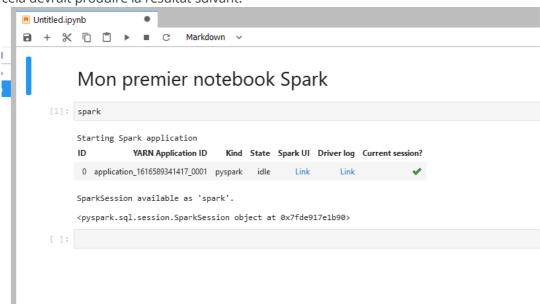
Spark : Spark en Scala

☐ SparkR : Spark en R.

Créez un notebook pyspark et exécutez dans un cellule la ligne suivante :

1 | spark

cela devrait produire la résultat suivant.



Pour des questions de sécurité, les liens vers l'interface Spark (*Spark UI*) et les journaux du driver (*Driver log*) ne fonctionnent pas. Il faut pour y accéder réaliser une connexion SSH + un transfert de port.

Exercice 3. First steps with Spark

3.1 Your first DataFrame — Spark's main object class is the DataFrame, which is a distributed table. It is analogous to R's or Python's data frames: one row represents an observation, one column represents a variable. But contrary to R or Python, Spark's DataFrames can be distributed over hundred of nodes.

• Run the following:

```
df = spark.read.format("parquet")
    .option("mode", "failFast")
    .option("header", "true")
    .option("inferSchema", True)
    .option("path", "file://path/to/file.csv.gz")
    .load()

df.cache()
```

You have just created a data frame! 🏂

Data frames are **immutable**: there is no method to alter one specific value once one is created. Also, data frames are **distributed**: they are split into blocks, ill-named **partitions** ¹, that are stored separately in the memory of the workers nodes.

The input file is a parquet file. Parquet is an open source column-oriented format that provide storage optimization. Spark natively can create a DataFrames from a parquet file.

Why do we cache the DataFrame? And are there any other solutions?

- In the Spark console, click on "Show incomplete applications" to see our current pyspark session. Look at the timeline. How many **executors** were used to perform the importation?
- Click on the *job* corresponding to the imporation, then on the unique *stage* that composes this job. On which nodes of the cluster were the executors located?
- For the importation stage, we have the equivalence one task = one partition, since the tasj is actually to create one partition. Does the number of tasks executed on each node recall you anything? (*Hint:* go back to the HDFS Console.)
- Lastly, open the event time line. From what you see, how many processors on each nodes were thre? Can you confirm it from EC2?

r emo::ji("wrench") **3.3 DataFrame manipulation** — Data frames are immutable, but they can be *transformed* in other data frames. Such *transformations* include: filtering, sampling, dropping columns, selecting columns, adding new columns...

• First, you can get information about the columns with:

```
flights.columns  # get the column names
flights.schema  # get the column names and their respective
type
flights.printSchema() # same, but human-readable
```

• What does the following code do?

```
passengers_per_month = flights\
select("PASSENGERS","YEAR","MONTH")\
groupBy("YEAR","MONTH")\
sum("PASSENGERS")
```

• And this one?

```
flights_from_2018 = flights\
sample(fraction=0.001)\
filter(flights.YEAR==2018)\
limit(100)
```

And this one?

```
overconfident_carriers = flights\
2
      .select("CARRIER", "DEPARTURES_SCHEDULED",
    "DEPARTURES_PERFORMED")\
3
     .withColumn(
                       # computes new variable
       "OVERCONFIDENCE", # called "OVERCONFIDENCE"
4
5
        (flights.DEPARTURES_SCHEDULED - flights.DEPARTURES_PERFORMED)/
6
       flights.DEPARTURES_PERFORMED
7
     .groupBy("CARRIER")\
8
9
      .sum("OVERCONFIDENCE")\
      .sort("sum(OVERCONFIDENCE)")
10
```

• Run each of the code sections.

r emo::ji("sleeping") 3.3 Lazy evaluation

- What happens when you run flights, like you would do in Python or R? Why?
- At question **3.2**, did you get any result at all? Did any of the instructions cause computation to actually happen? (*Hint:* look at the Spark console)

This is because Spark has what is known as **lazy evaluation**, in the sense that it will wait as much as it can before performing the actual computation. Said otherwise, when you run an instruction such as:

```
1 | filtered_flights = flights.filter(fligths.YEAR==2018)
```

... you are not executing anything! Rather, you are building an **execution plan**, to be realised later.

Spark is quite extreme in its lazyness, since only a handful of methods called **actions**, by opposition to **transformations**, will trigger an execution. The most notable are:

- 1. collect(), explicitly asking Spark to fetch the resulting rows instead of to lazily wait for more instructions,
- 2. take(n), asking for n first rows
- 3. first(), an alias for take(1)
- 4. show() and show(n), human-friendly alternatives ²
- 5. count(), asking for the numbers of rows
- 6. all the "write" methods (write on file, write to database)

This has advantages: on huge data, you don't want to accidently perform a computation that is not needed. Also, Spark can optimize each **stage** of the execution in regard to what comes next. For instance, filters will be executed as early as possible, since it diminishes the number of rows on which to perform later operations. On the contrary, joins are very computation-intense and will be executed as late as possible. The resulting **execution plan** consists in a **directed acyclic graph** (DAG) that contains the tree of all required actions for a specific computation, ordered in the most effective fasshion.

This has also drawbacks. Since the computation is optimized for the end result, the intermediate stages are discarded by default. For instance, in the following:

```
1  # step 1
   flights_overconfidence = flights\
3
     .withColumn(
4
        "OVERCONFIDENCE",
5
        (flights.DEPARTURES_SCHEDULED - flights.DEPARTURES_PERFORMED)/
        flights.DEPARTURES_PERFORMED
6
7
     )
8 # step 2
9
   flights_overconfidence_2018 = flights_overconfidence\
     .filter(fligths.YEAR==2018)
10
11
      .collect()
```

... the intermediate flights_overconfidence does not exist more after collect() have been called than before the call. Indeed, the values for other years than 2018 have not be computed at all!

Now run:

```
passengers_per_month.show()
flights_from_2018.count()
overconfident_carriers.take(10)
```

Was something executed this time?

• You can get the execution plan from the Spark console, or from Python with the explain() method. Try with flights_from_2018.explain(). Does the order of the stages make sense?

```
r emo::ji("billed_hat") 3.4 Practice
```

The complete list of methods (transformations and actions) for data frames is listed here. The aggregation functions, such as sum(), max(), mean()... are listed here.

- What are the 10 biggest airports of the USA in 2018?
- What is the longest regular flight served by each carrier in 2000?

Exercice 4. Map-and-reduce architecture

Manny computation algorithms can be expressed using two stages:

• a **map stage**, where the intructions can be applied element-wise, in the sense the if elements are arranged as a list, the operation on element e, does not depend on the value of e^\prime

• a **reduce stage**, where the instructions obtained in the first stage are combined pairwise recursively; each time a result is obtained from the first stage, it is combined with earlier results, as in an accumulator

The reduce function must be associative, and commutativity simplifies the reduce step even further. Typical exemples are addition and multiplication. Concatenation is associative, but not commutative.

r emo::ji("man_teacher") 4.1. Map-and-reduce exemples

- Find two exemples of computation problems that decompose well under the map-and-reduce principle, and one that can't.
- The count() method is expressible as a *map-and-reduce* algorithm. flights.count() is equivalent to the following code. Can you make clear how the job is executed? Is it faster?

```
1 # the map function is not available at the data frame level
2 | # we have to go down at the data set (RDD) level
3 flights\
4
    .rdd\
5
    .map(lambda flight: 1)\
    .reduce(
6
7
      lambda accumulator, value:
8
         accumulator + value
9
     )
10 | # reduce is an action verb
# we do not need an explicit collect()
```

- Explain the lambda flight: 1 syntax. How do you call this kind of object?
- Look at the Spark console to see where the different stages of the computation actually happenned.

```
r emo::ji("soccer") 4.2. Practice
```

- Compute the total number of passengers transported following the map-and-reduce principle. Is it faster than
- What does the following code do?

```
def my_function( a, b ) :
    return b if b > a else a

flights\
    .rdd\
    .map(lambda flight: flight.AIR_TIME)\
    .reduce( my_function )
```

• The *map* stage may as well return a tupple (FR: n-uplet), as long as the you have an corresponding well chosen *reduce* stage. For instance, what does the following do?

```
flights\
    .rdd\
    .map(lambda flight: (flight.AIR_TIME, flight.CARRIER))\
    .reduce(lambda a, b: a if a[0] > b[0] else b)
```

- How would you recode the mean() function in two succesive map-and-reduce operations? Is it possible with only one?
- What about the variance?

^{1.} Usually a "partition" is an set of chunks that cover all the data, without any repetition between the chunks. But not in Spark! $\underline{\underline{e}}$

^{2.} first() is exactly take(1) (\underline{ref}) and show prints the result instead of returning it as a list of rows (\underline{ref}) \underline{e}